# ESTIMATING EVOLUTIONARY RATES AND DIVERGENCE TIMES

## Philippe Lemey[1], Guy Baele[1] and Marc Suchard[2]

[1] Rega Institute, Department of Microbiology and Immunology, K.U. Leuven, Belgium.

[2] Departments of Biomathematics and Human Genetics, David Geffen School of Medicine at UCLA. Department of Biostatistics, UCLA School of Public Health

---

**MOLECULAR SEQUENCES**

*Alignment Methods*                    BIOINFORMATICS

↓

**ALIGNMENT**

*Sequence Evolution Models*
*Phylogenetic Methods*                 PHYLOGENETICS

↓

**EVOLUTIONARY TREE**
(time scale = genetic distance)

*Molecular Clock Models*               PHYLOGENETICS

↓

**EVOLUTIONARY TREE**
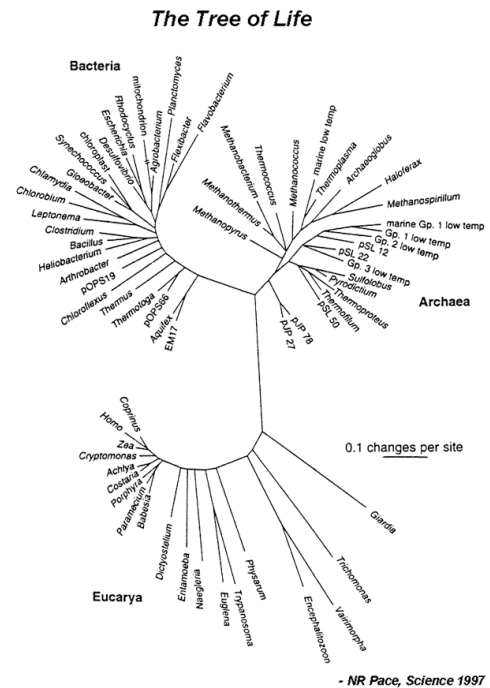(time scale = years)

*Coalescent Models*                    POPULATION GENETICS

↓

**EPIDEMIOLOGY**

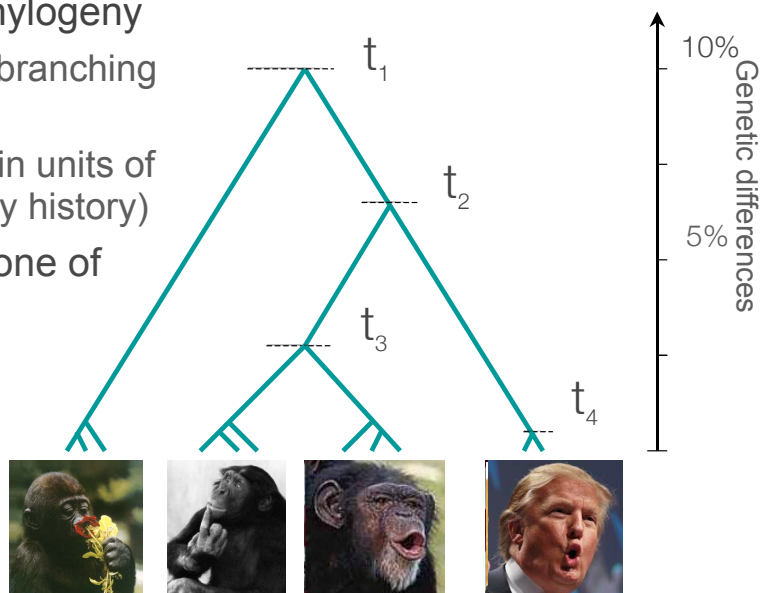# Molecular phylogenies

◉ most molecular phylogenies

‣ are unrooted (or the rooting is due to prior information)

‣ have branch lengths representing genetic change



*The Tree of Life*

0.1 changes per site

- NR Pace, Science 1997

---

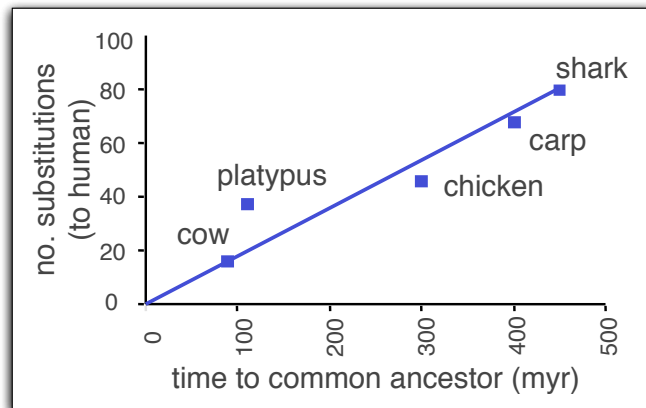# Molecular phylogenies

◉ the ideal molecular phylogeny
‣ is rooted (implies a branching order)
‣ has branch lengths in units of time (an evolutionary history)

◉ how do we construct one of these trees?



$t_1$

$t_2$

$t_3$

$t_4$

10%

5%

Genetic differences
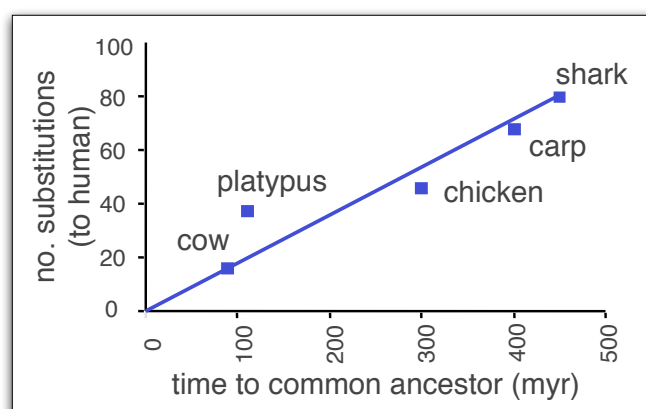
# A constant evolutionary rate through time

- to obtain a time phylogeny, the evolutionary model must assume a relationship between the accumulation of genetic diversity and time



- Zuckerkandl and Pauling (1962): the rate of amino acid replacements in animal haemoglobins was roughly proportional to real time, as judged against the fossil record

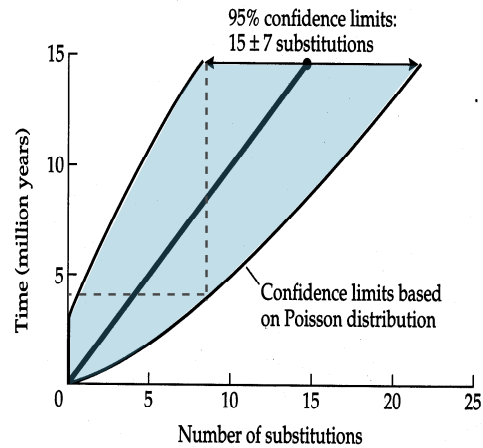# A constant evolutionary rate through time

- the *molecular clock* is particularly striking when compared to the obvious differences in rates of morphological evolution...

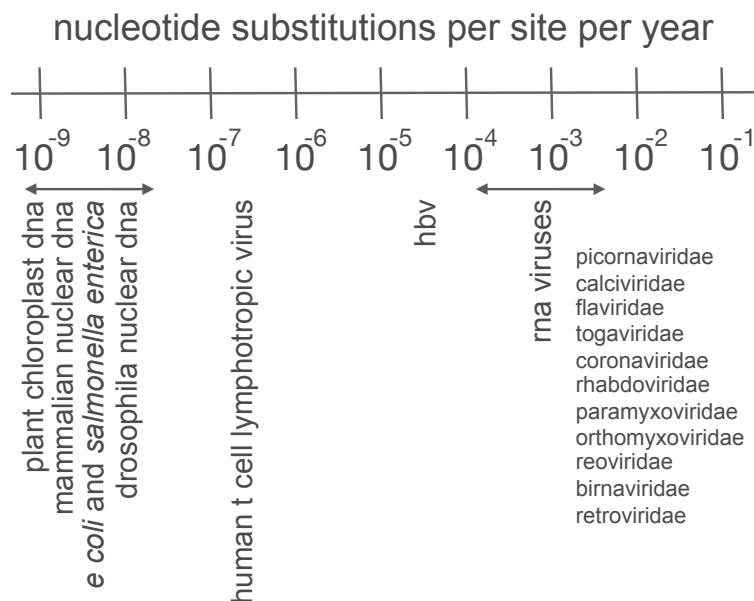# The molecular clock is not a metronome

- if mutation every MY with Poisson variance

  ‣ 95% of the lineages 15MY old have 8-22 substitutions
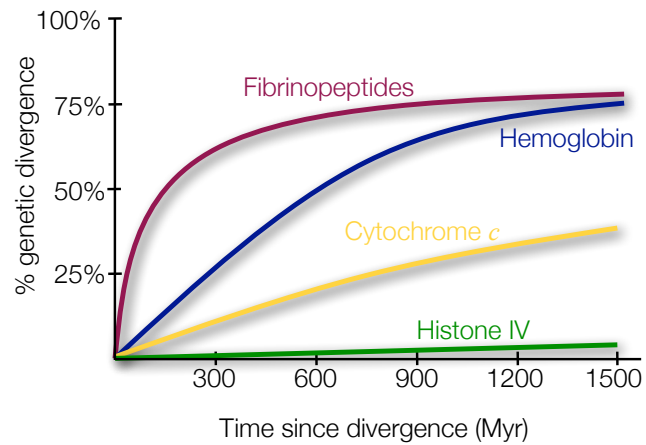
  ‣ 8 substitutions also could be < 5 MY old



‣ Molecular Systematics, p532.

---

# And there is no global molecular clock

nucleotide substitutions per site per year

# And there is no global molecular clock

- different genes, different profiles
- variation in mutation rate?
- variation in selection

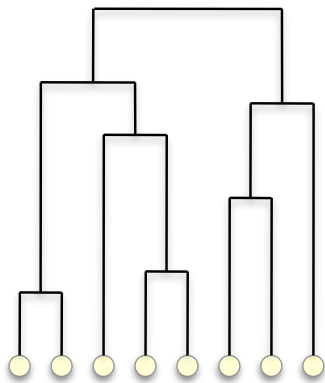  genes coding for some molecules under very strong stabilizing selection



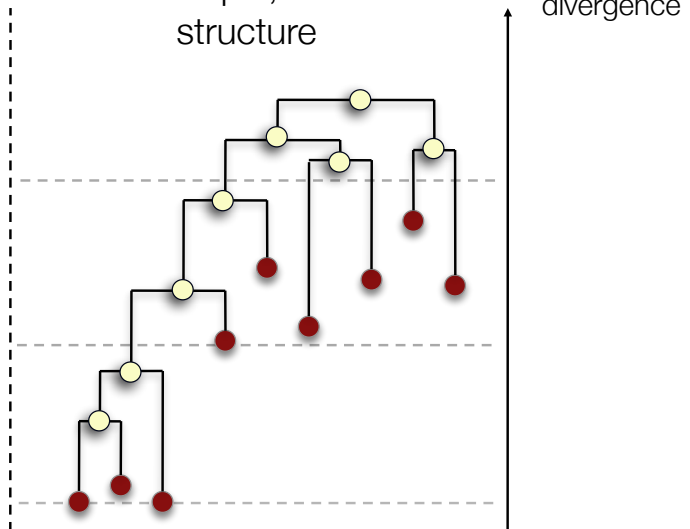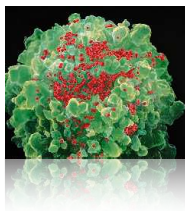---

calibrating the molecular clock

# From substitution units to time units



nodes with
point calibrations

22 Mya

7 Mya

Contemporary sample
probabilistic calibrations

time

95% C.I.
20-30 Mya

95% C.I.
5-15 Mya

now

# Node Calibrations



**Fossils**

**biogeography**

Kauai
(5.1 My)

Oahu
(3.7 My)

Maui-Nui (W. Molokai)
(1.9/1.6 My)

0    60    120    km

0.01

Hawaii
(0.43 My)

Main Hawaiian Islands
(K-Ar ages)

**host-pathogen co-divergence**

| | |
|---|---|
| Felis catus | FdPV1 |
| Puma concolor | PcPV1 |
| Lynx rufus | LrPV1 |
| Panthera leo | PlpPV1 |
| Panthera uncia | UuPV1 |
| Pli | PlPV1 |
| Cfe | COPV |

100   100   100   95   100   100   89   83   93   100   100   100   100   100

0.01    0.1

# Calibration using sampling times



contemporary sample,
no time structure

serial sample, with time
structure

divergence

# Tip calibration: two major applications



RNA viruses
evolve quickly:
$10^{-3}$ - $10^{-5}$
substitutions per
site per year.

- Substitutions accumulate
  between the times of sampling

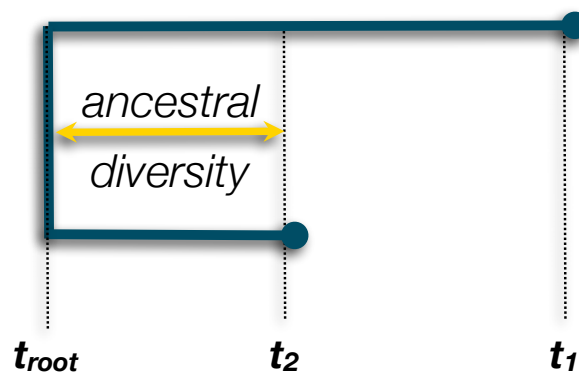- Serially sampled sequences or
  heterochronous sequences

*Measurably evolving population*



ancient DNA
data sets of
radiocarbon-dated
specimens

# incorporating sampling time: naive method

observed number of substitutions
or genetic divergence
d

sampling time 1
$t_1$

sampling time 2
$t_2$

substitution rate, $\mu$
$= d / |t_1 - t_2|$



# incorporating sampling time: naive method

*ancestral*

*diversity*

$t_{root}$         $t_2$         $t_1$

# incorporating sampling time: naive method



$$\mu = (d_1 - d_2) / (t_1 - t_2)$$

# linear regression



$$\mu = d_i / (t_i - t_{root})$$

- can be rearranged:

$$d_i = \mu \, (t_i - t_{root})$$

$$E[d_i] = \mu \cdot t_i - \mu \cdot t_{root}$$

gradient is: $\mu$

y-intercept is: $-\mu \cdot t_{root}$

x-intercept is: $t_{root}$

# Estimating the time-scale

- H1N1/09 'Swine Flu'
- Rate: $3.14E^{-3}$
  mutations/genomic site/year
- tMRCA: 2009.041
  (15-Jan-2009)
- Correlation: 0.83
- $R^2$: 0.69



# A DNA virus (smallpox)

Variola, *Poxviridae*, 190kb genome
Sampling 1946-1977

**VARV**
$R^2=0.67862$

Rate estimate: $8.2 \times 10^{-6}$ Subs/Site/Year

# Salmonella Typhimurium



# Exploring the temporal structure of heterochronous sequences using TempEst (formerly Path-O-Gen) 🔓 ⓒ

Andrew Rambaut, Tommy T. Lam, Luiz Max Carvalho, Oliver G. Pybus

# Two lost decades of seasonal H1N1 evolution



# Two lost decades of seasonal H1N1 evolution



A

B

20 years

0.05 subst/site/year

Root-to-tip distance

- Linear regression of genetic distance against sampling time shows no divergence between 1957 and 1977 (when H1N1 re-emerged).

- Also apparent when comparing molecular clock trees (A) with non-clock trees (B).

- ght ab

*1*

## Time structure via tip calibration



Contemporary sample
no time structure

Serial sample
with time structure

time

1980

1990

2000

‣ Rambaut A. (2000) *Bioinformatics*, **16**, 395-399.

---

## Clock versus non-clock

- strict molecular clock:
  Zuckerkandl & Pauling (1962) in Horizons in Biochemistry, pp. 189–225
  - ‣ all lineages evolve at the same rate
  - ‣ allows the estimation of the root of the tree and dates of individual nodes
- unconstrained (unrooted) Felsenstein model:
  Felsenstein (1981) *JME*, **17**: 368 - 376
  - ‣ each branch has its own rate independent of all others
  - ‣ time and rate are confounded and can only be estimated as a compound parameter (branch lengths)

## Likelihood ratio test for molecular clock models

- complex model $H_1$



- null model $H_0$



2N-3 parameters                    N-1 parameters
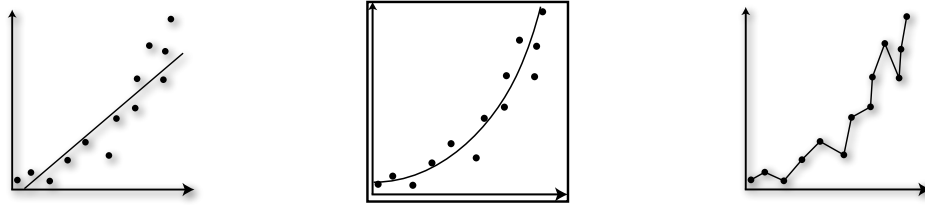
$LRS = 2(\max[\ln L(H_a|D)] - \max[\ln L(H_0|D)])$

- likelihood ratio test with N-2 degrees of freedom
- models are nested because values of $b_1$-$b_7$ can be specified that give node heights $t_1$-$t_4$

---

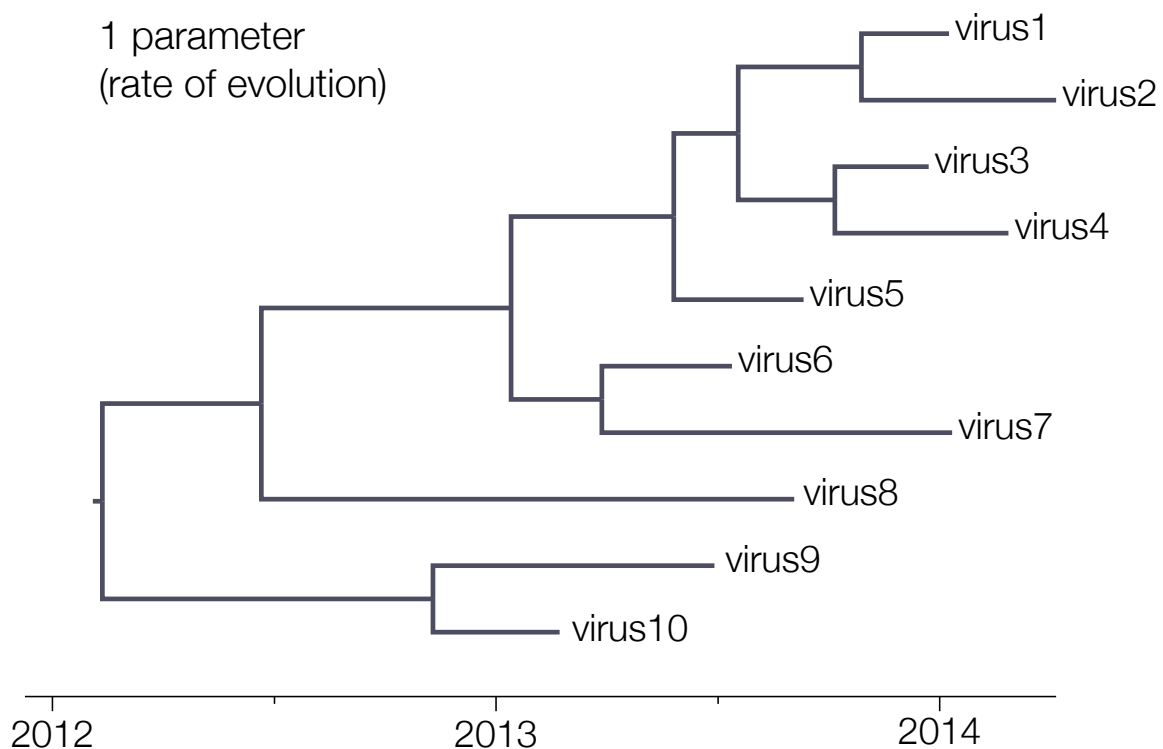## Relaxing the molecular clock

# Need for a relaxed molecular clock

- the unrooted model of phylogeny and the strict molecular clock model are two extremes of a continuum.
- dominate phylogenetic inference
- but both are biologically unrealistic:
  - ‣ the real evolutionary process lies between these two extremes
  - ‣ model misspecification can produce positively misleading results



‣ Pybus (2006) *Genome Biol*. **4**, e151
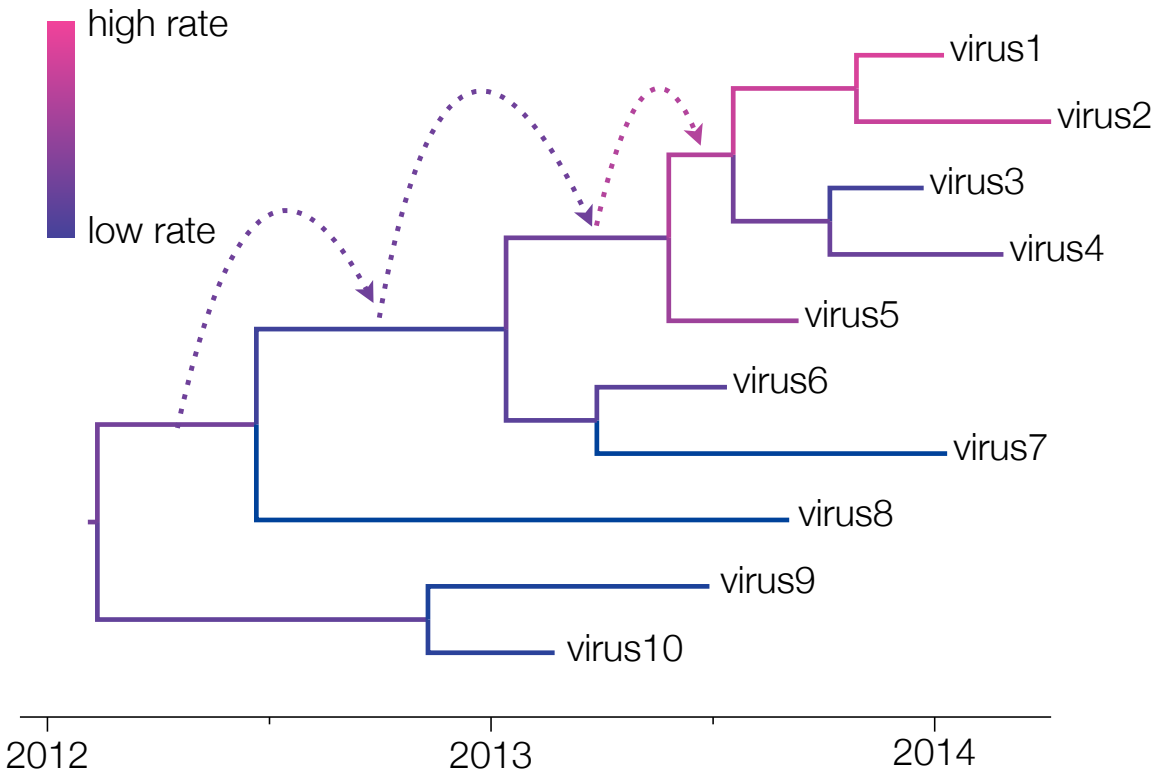
---

# 'strict' molecular clock

1 parameter
(rate of evolution)



| | | | |
|---|---|---|---|
| 2012 | 2013 | 2014 | |

'local' molecular clock

high rate

low rate

virus1
virus2
virus3
virus4
virus5
virus6
virus7
virus8
virus9
virus10

2012          2013          2014

host specific local clock

high rate

low rate

pig▸human

bird▸pig

human
human
human
human
pig
pig
pig
bird
bird
bird

2012          2013          2014

autocorrelated relaxed clock

high rate
low rate

virus1
virus2
virus3
virus4
virus5
virus6
virus7
virus8
virus9
virus10

2012    2013    2014

lognormal uncorrelated relaxed clock

low rate    high rate

virus1
virus2
virus3
virus4
virus5
virus6
virus7
virus8
virus9
virus10

2 parameters
(mean rate and
variance in rate among
branches)

2012    2013    2014

# Relaxed clocks: (1) local molecular clocks



‣ specify $H_0$ beforehand

‣ problem of identifiability

‣ Yoder and Yang (2000) Mol Biol & Evol **17**: 1081-1090.

---

# Bayesian local clocks



*Worobey et al., Nature, 2014; 508(7495): 254–257*

# Autocorrelated relaxed clocks

- rates for each branch are drawn from a distribution centered on the rate of the ancestor
  - ‣ but what is the rate at the root?
  - ‣ A prior degree of autocorrelation?
  - ‣ not currently possible to do phylogenetic inference



$$r_i \sim LogNormal(r_{A(i)}, \sigma^2 \Delta t_i)$$

‣ e.g., Thorne JL, Kishino H, Painter IS (1998) Mol Biol & Evol **15**: 1647-1657.

---

# Uncorrelated relaxed clocks

- rates for each branch are drawn independently from an identical distribution:



$$r \sim Exp(\lambda)$$

$$r \sim LogNormal(\mu, \sigma^2)$$

$$r \sim Gamma(\alpha, \beta)$$

‣ Drummond et al. (2006) Plos Biology **4**: e88.

# Bayesian evolutionary analysis sampling trees

- Given sequence data that is temporally spaced estimate true values of:
  - ‣ substitution parameters ($\mu$ and $Q$)
  - ‣ ancestral genealogy ($g = E_g$, $t_Y$)
    - tree topology
    - dates of divergence
  - ‣ population history ($\theta$)



$\mu$     $Q$

- Bayesian inference

$$P(g,\mu,\theta,Q|D) = \frac{1}{Z} Pr\{D|g,\mu,Q\} f_g(g|\theta) f_\mu(\mu) f_\theta(\theta) f_Q(Q)$$
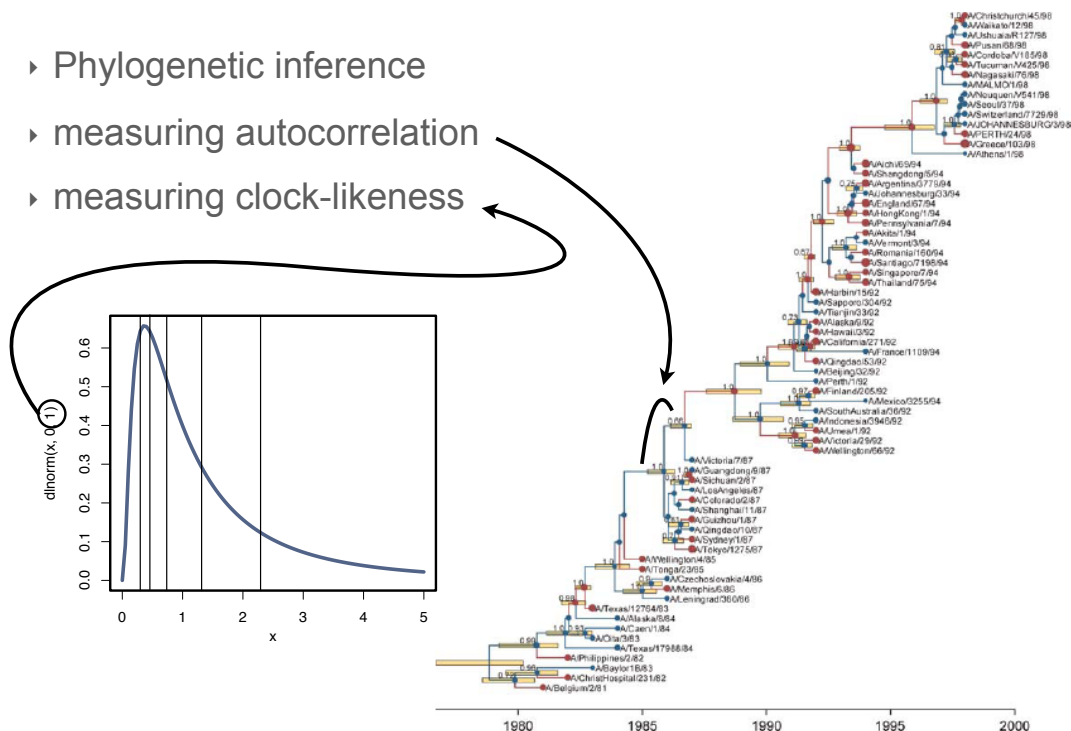
"relaxed phylogenetics and dating with confidence"

$t = \{ t_1, t_2, \ldots, t_{2n-1} \}$

$R = \{ r_1, r_2, \ldots, r_{2n-1} \}$     $f(R|g) = f(R) = \prod_{i=1} \lambda e^{-\lambda r_i}$



$N_e$

time

---

# Uncorrelated relaxed clocks: example

- ‣ Phylogenetic inference
- ‣ measuring autocorrelation
- ‣ measuring clock-likeness

## Evaluating clock-like behaviour?



## Model testing using Bayes factors

- Bayes factor $\qquad B_{01} = \dfrac{p(Y|M_1)}{p(Y|M_0)}$

- when two models $M_0$ and $M_1$ are being compared, one defines the Bayes factor in favor of $M_1$ over $M_0$ as the ratio of their respective marginal likelihoods

- When there are unknown parameters, the Bayes Factor $B_{01}$ has in a sense the form of a *likelihood ratio*
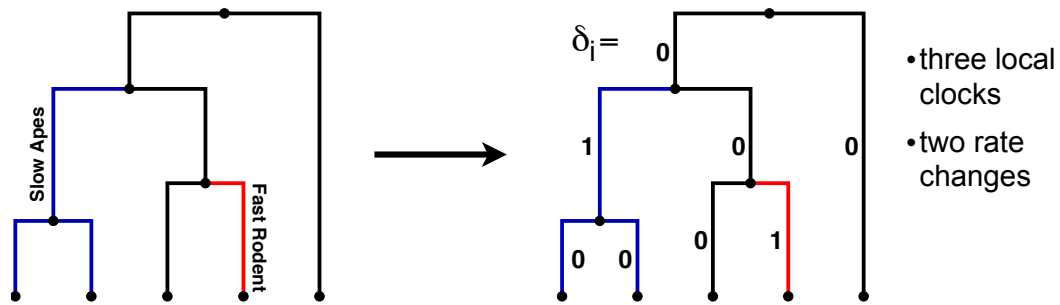
Guy

# Random local clocks

➡ critics on the local clocks
- specify H0 a priori
- problem of identifiability

➡ critics on the uncorrelated relaxed clocks
- Rate changes do not necessarily occur regularly or on every branch
- Small number of significant changes

*So, can we handle the uncertainty in the number and locations of a small number of local clocks?*



$\delta_i =$

- three local clocks
- two rate changes

➡ How to explore $2^{2n-2}$ clock models?

---

# Random local clocks

➡ Using Bayesian stochastic search variable selection: formulate a prior that such that many rate changes (indicators) are 0 but allow the data to determine which ones are required to explain (most of the) rate variation using MCMC



➡ Three mtDNA nuclear genes from 42 mammals (Douzery, 2003)

➡ 5-12 local clocks

*Drummond and Suchard, 2010.*

# Random local clocks

➡ Testing whether a branch accommodates a rate change using Bayes factors

◉ Data D is assumed to have been arisen under one of two models, or one of two hypotheses $H_1$ and $H_2$.

$$pr(H_k|\mathbf{D}) = \frac{pr(\mathbf{D}|H_k)pr(H_k)}{pr(\mathbf{D}|H_1)pr(H_1) + pr(\mathbf{D}|H_2)pr(H_2)}$$

so that

$$\boxed{\frac{pr(H_1|\mathbf{D})}{pr(H_2|\mathbf{D})}} = \boxed{\frac{pr(\mathbf{D}|H_1)}{pr(\mathbf{D}|H_2)}}\boxed{\frac{pr(H_1)}{pr(H_2)}}$$

*posterior odds*   *Bayes factor*   *prior odds*



◉ Prior probabilities pr($H_1$) and pr($H_2$) = 1 - pr($H_1$). Posterior probabilities pr($H_1$|D) and pr($H_2$|D) = 1 - pr($H_1$|D)

---

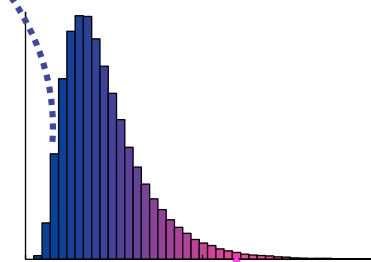# Extensions for testing evolutionary rate hypotheses

Pybus and Rambaut, NGR, 2009

**New insights into the evolutionary rate
of HIV-1 at the within-host and
epidemiological levels**

Katrina A. Lythgoe* and Christophe Fraser

months

years

Lemey et al 2006 AIDS Rev

**New insights into the evolutionary rate
of HIV-1 at the within-host and
epidemiological levels**

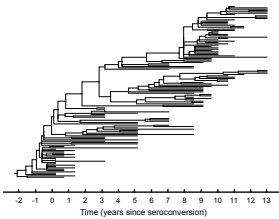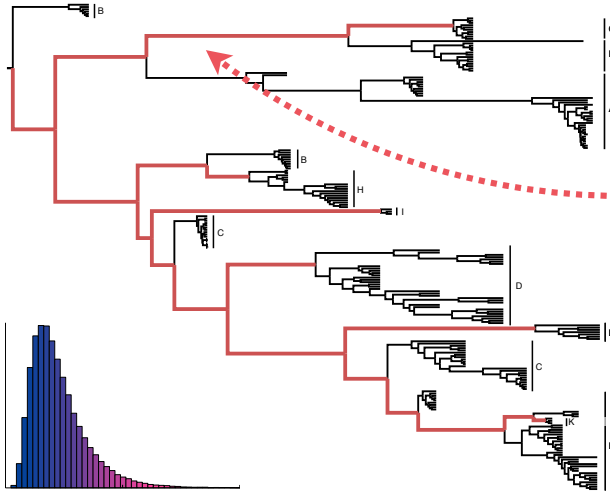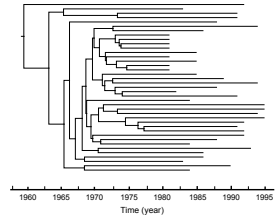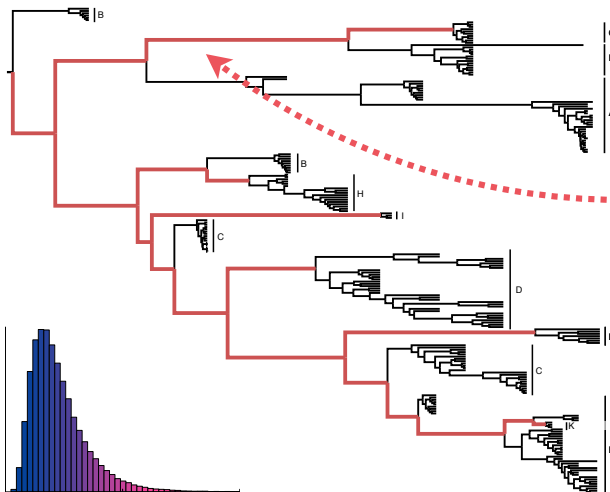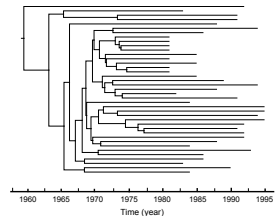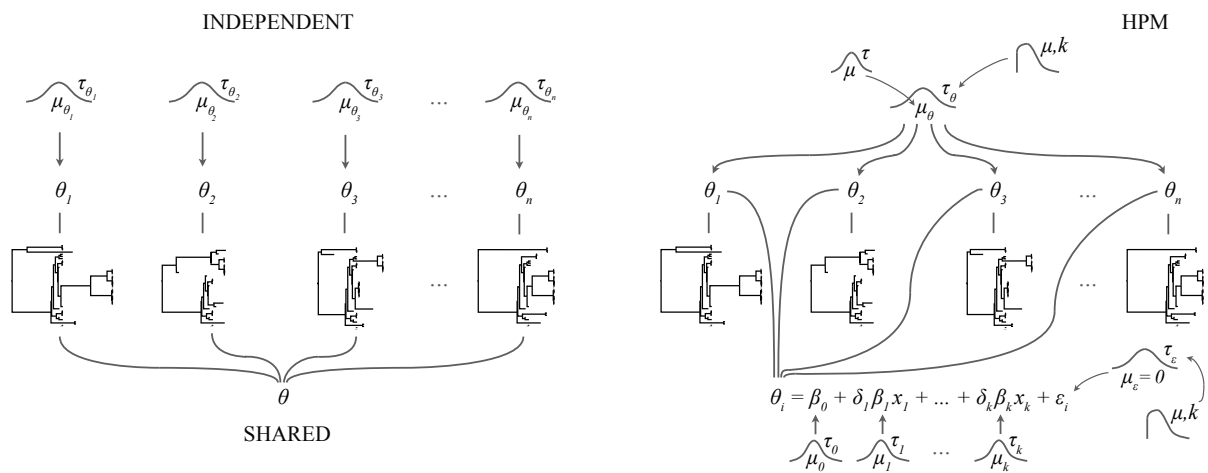Katrina A. Lythgoe* and Christophe Fraser

Time (years since seroconversion)

Time (year)

random effects model:

$$\log \mu_i = \theta_i$$

low rate    high rate

Vrancken et al., PLoS Comp Bio, 2014

mixed effects model:

$$\log \mu_i = \theta_i + \beta X_i$$
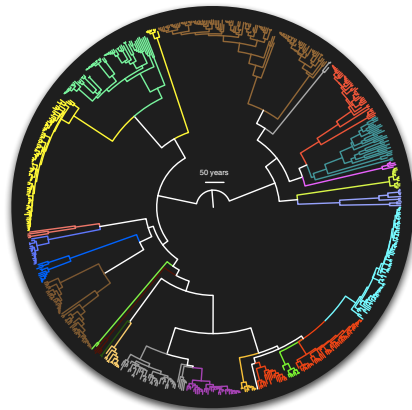
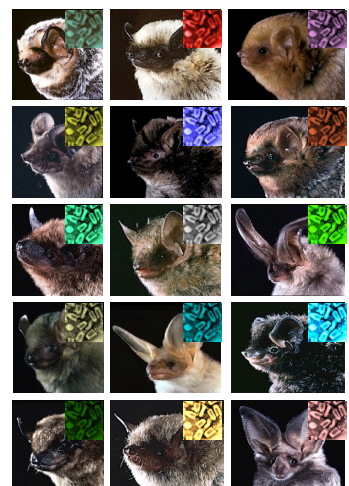| | pol | env |
|---|---|---|
| Rate ($10^{-3}$ subst./site/yr) | | |
| $X_i$=0 (within host) | 5.70 (4.02-6.21) | 10.37 (8.06-12.76) |
| $X_i$=1 (transmitted lineage) | 2.21 (1.57-2.99) | 3.80 (2.32-5.20) |
| ln Bayes factor ($rate_{transmitted} < rate_{within}$) | | |
| | >7.50 | >6.29 |

Vrancken et al., PLoS Comp Bio, 2014

# Hierarchical phylogenetic modelling

INDEPENDENT

HPM

$$\theta_i = \beta_0 + \delta_1 \beta_1 x_1 + \ldots + \delta_k \beta_k x_k + \varepsilon_i$$
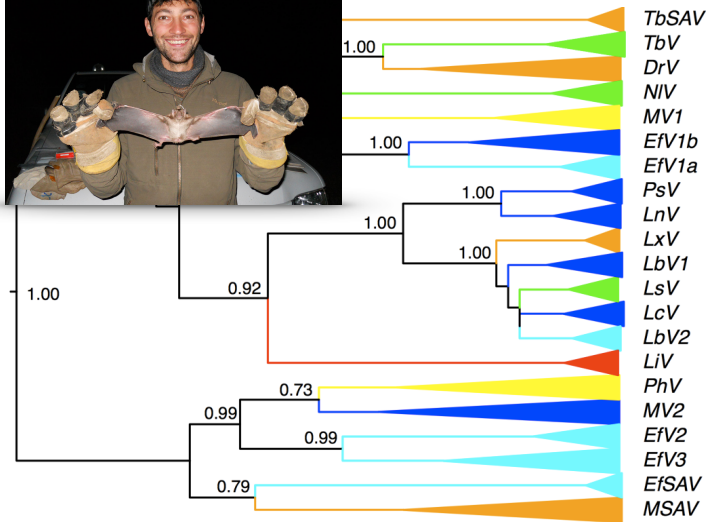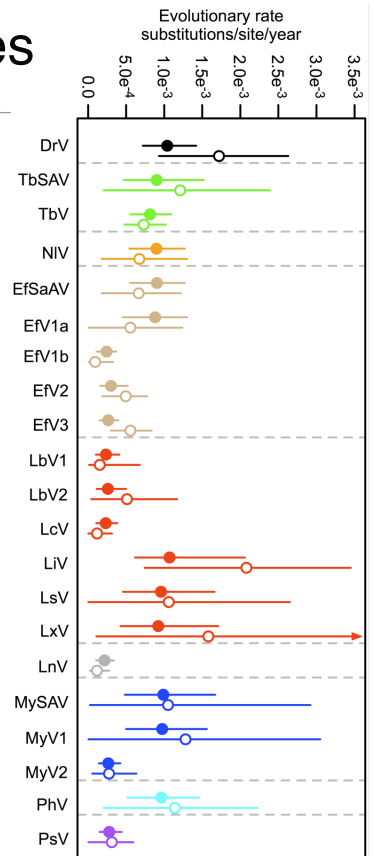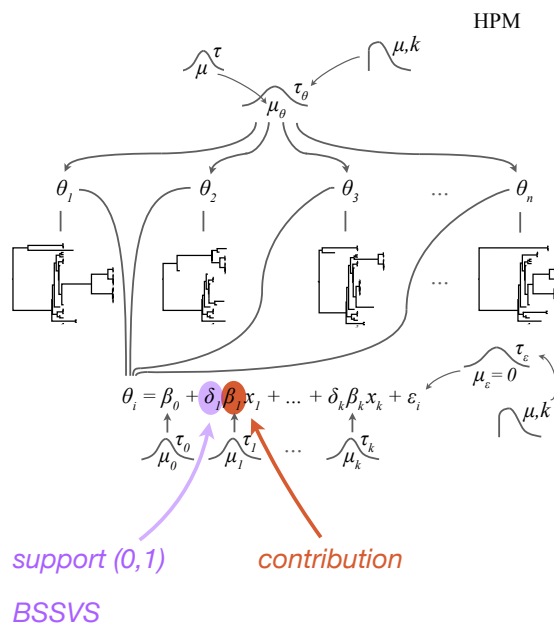
SHARED

*…with fixed effects*



*Courtesy of D. Streicker*
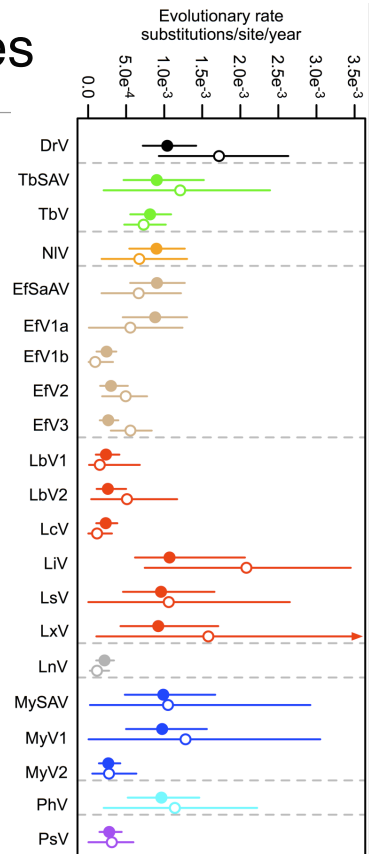
# Bat rabies virus evolutionary rates

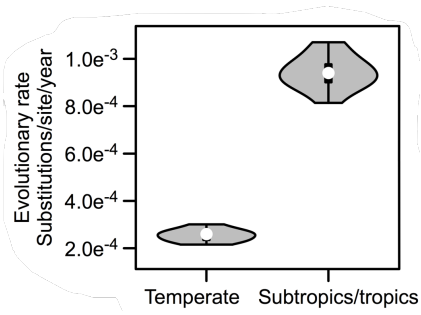Streicker et al., 2012. *PLoS Pathogens*

# Bat rabies virus evolutionary rates

HPM

$$\theta_i = \beta_0 + \delta_1\beta_1 x_1 + ... + \delta_k\beta_k x_k + \varepsilon_i$$

*support (0,1)*     *contribution*

*BSSVS*

Streicker et al., 2012. *PLoS Pathogens*

# Bat rabies virus evolutionary rates



| Predictor | Bayes factor | $\beta$ (95% HPD) \| $\delta = 1$ |
|---|---|---|
| Climate | 466.54 | |
| Basal metabolic rate | 0.82 | |
| Torpid metabolic rate | 1.00 | |
| Coloniality | 0.46 | |
| Seasonal activity | 0.46 | |
| Long-distance migration | 0.69 | |

Streicker et al., 2012. *PLoS Pathogens*